

CONSIDERING MODEL UNCERTAINTY IN RANDOM REGRESSION MODELS BY MEANS OF BAYESIAN VARIABLE SELECTION

J.P. Steibel¹ and F. Grignola²

¹Facultad de Agronomía, Universidad de Buenos Aires, Cap. Fed., Argentina

²Monsanto Company, Saint Louis, MO, USA

INTRODUCTION

In animal breeding, many traits admit repeated measurements or test-day records over time. In order to model the expectation and covariance of these characters as a function of time, the interest in using test-day models (TDM) in the context of the Gaussian mixed linear model has increased in recent years. Even though the advantages of TDM are well known, a problem that frequently arises is the choice of suitable number terms in the linear function. Jensen (2001) discussed various strategies about model choice and, within the Bayesian framework, he suggested Bayesian model averaging to consider the uncertainty around the model for the prediction of breeding values. In this work we implement a simple and flexible Bayesian approach proposed by Kuo and Mallick (1998) to subset a pre-specified set of covariates that best describe the trait of interest in a random coefficient regression model (RRM). The posterior probability of each regressor entering the model is computed using the Gibbs sampling algorithm. The method is illustrated with a simple example where the variable selection strategy is limited to the fixed effects.

METHODS

The method of Kuo and Mallick (1998) expands the usual regression model to include an indicator variable for each predictor considered. In standard notation and for a RRM,

$$y_{jkl} = \sum_{m=0}^{nb} b_{km} \gamma_{km} z_{jlm} + \sum_{m=0}^{nu} u_{jm} \delta_{km} z_{jlm} + \sum_{m=0}^{nu} p_{jm} \theta_{km} z_{jlm} + e_{jkl}$$

where y_{jkl} is the record l for animal j in contemporary group k , z_{jlm} is the known covariate m for animal j and record l associated with time, b_{km} are the “fixed” regression coefficients that describe the average curve for contemporary group (CG) k , u_{jm} and p_{jm} the additive genetic and permanent environmental random regression coefficients for animal j , and γ , δ and θ are indicator variables (I). In matrix notation

$$\mathbf{Y} = \mathbf{X}\Gamma\mathbf{b} + \mathbf{Z}\Delta\mathbf{u} + \mathbf{W}\Theta\mathbf{p}, \text{ with } \begin{pmatrix} \mathbf{b} \\ \mathbf{u} \\ \mathbf{p} \\ \mathbf{e} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{b}_0 \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}; \begin{bmatrix} \mathbf{D}_0 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G} \otimes \mathbf{A} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{P} \otimes \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{R} \end{bmatrix} \right)$$

where \mathbf{D}_0 , \mathbf{G} and \mathbf{P} are the covariance matrices for the regressors describing the “fixed”, additive genetic and permanent environmental effects, \mathbf{A} = additive relationship matrix, $\Gamma = \text{Diag}(\gamma_{km})$, $\Delta = \text{Diag}(\delta_{jm})$, $\Theta = \text{Diag}(\theta_{jm})$ and $\mathbf{R} = I\sigma_e^2$.

Prior distributions are : $\mathbf{G} \sim \text{Inv-Wishart}_{\text{nu}}((v_a, (v_a \mathbf{S}_a)^{-1})$; $\mathbf{P} \sim \text{Inv-Wishart}_{\text{nu}}((v_p, (v_p \mathbf{S}_p)^{-1})$; $\sigma_e^2 \sim \text{I-Gamma}(v_e, s_e)$; γ, δ and θ are i.i.d. with Binomial distribution $\text{Bi}(1, p_i)$, where p_i = prior probability of including the i^{th} regression coefficient in the model. The joint posterior density is $f(\Gamma, \beta, \Delta, \Theta, p, \sigma_a^2, \sigma_p^2, \sigma_e^2 | \mathbf{y}) \propto f(\mathbf{y} | \Gamma, \beta, \sigma_e^2) f(\Gamma) f(\beta) f(\sigma_e^2) f(\Delta) f(\mathbf{u} | \mathbf{G}) f(\mathbf{G}) f(\Theta) f(\mathbf{p} | \mathbf{P}) f(\mathbf{P})$, and the posterior conditional distributions are

$$\begin{pmatrix} \mathbf{b} \\ \mathbf{u} | \mathbf{Y}, \mathbf{G}, \mathbf{P}, \mathbf{R}, \Gamma, \Delta, \Theta \\ \mathbf{p} \end{pmatrix} \sim N(\mathbf{C}^{-1} \mathbf{r}; \mathbf{C}^{-1}), \mathbf{r} = \begin{pmatrix} \Gamma' \mathbf{X}' \mathbf{R}^{-1} \mathbf{Y} + \mathbf{D}_\theta^{-1} \mathbf{b}_0 \\ \Delta' \mathbf{Z}' \mathbf{R}^{-1} \mathbf{Y} \\ \Theta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{Y} \end{pmatrix}$$

$$\mathbf{C} = \begin{pmatrix} \Gamma' \mathbf{X}' \mathbf{R}^{-1} \mathbf{X} \Gamma + \mathbf{D}_\theta^{-1} & \Gamma' \mathbf{X}' \mathbf{R}^{-1} \mathbf{Z} \Delta & \Gamma' \mathbf{X}' \mathbf{R}^{-1} \mathbf{W} \Theta \\ \Delta' \mathbf{Z}' \mathbf{R}^{-1} \mathbf{X} \Gamma & \Delta' \mathbf{Z}' \mathbf{R}^{-1} \mathbf{Z} \Delta + \mathbf{G}^{-1} \otimes \mathbf{A}^{-1} & \Delta' \mathbf{Z}' \mathbf{R}^{-1} \mathbf{W} \Theta \\ \Theta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{X} \Gamma & \Theta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{Z} \Delta & \Theta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{W} \Theta + \mathbf{G}^{-1} \otimes \mathbf{A}^{-1} \end{pmatrix}$$

$$\begin{aligned} [\mathbf{G} | \mathbf{Y}, \mathbf{b}, \mathbf{u}, \mathbf{p}, \mathbf{P}, \mathbf{R}, \Gamma, \Delta, \Theta] &= IW(q_g + v_a; (\mathbf{u}' \mathbf{A}^{-1} \mathbf{u} + v_a \mathbf{S}_a)^{-1}) \\ [\mathbf{P} | \mathbf{Y}, \mathbf{b}, \mathbf{u}, \mathbf{p}, \mathbf{G}, \mathbf{R}, \Gamma, \Delta, \Theta] &= IW(q_p + v_p; (\mathbf{p}' \mathbf{p} + v_p \mathbf{S}_p)^{-1}) \\ [\sigma_e^2 | \mathbf{Y}, \mathbf{b}, \mathbf{u}, \mathbf{p}, \mathbf{G}, \mathbf{P}, \mathbf{R}, \Gamma, \Delta, \Theta] &= IG(n + v_e; \hat{\rho}' \hat{\rho} + v_e \cdot S_e), \text{ with } \hat{\rho} = \mathbf{Y} - \mathbf{X} \Delta \mathbf{b} - \mathbf{Z} \Delta \mathbf{u} - \mathbf{W} \Theta \mathbf{p}. \end{aligned}$$

The posterior distribution of each indicator variable (I_i) is : $(I_i | \mathbf{Y}, \mathbf{b}, \mathbf{u}, \mathbf{p}, \mathbf{G}, \mathbf{P}, \mathbf{R}, I_i) \sim \text{Bi}(1, \omega_i)$

$$\omega_i = \frac{p_i \cdot \exp(\mathbf{e}^* \mathbf{R}^{-1} \mathbf{e}^*)}{(1 - p_i) \cdot \exp(\mathbf{e}^{**} \mathbf{R}^{-1} \mathbf{e}^{**}) + p_i \cdot \exp(\mathbf{e}^* \mathbf{R}^{-1} \mathbf{e}^*)}, \mathbf{e}^* = \text{residual if } I_i=1 \text{ and } \mathbf{e}^{**} = \text{residual if } I_i=0$$

The method allows to select a different linear function for each effect in the model or to force the same function in some of them.

Linear growth model example. We present a simplified example where the uncertainty around the model is limited to the “fixed” effects. The data came from the beef cattle selection experiment included in the DFREML package (Meyer, 1998). The model was defined as

$$y_{jkl} = \sum_{m=0}^5 \left(\prod_{i=0}^m (y_{ki}) b_{km} z_{jlm} \right) + \sum_{m=0}^5 u_{jm} z_{jlm} + \sum_{m=0}^5 p_{jm} z_{jlm} + e_{jkl}$$

Note that in this model the specification of the indicator variables is such that it is only possible to select complete polynomials. If one $I_i = 0$, regressors of higher order will be excluded. Legendre polynomials were used to model the “fixed” and random curves, where m determines the number of polynomial terms in the model. The contemporary groups represented 10 paddocks and the maximum m allowed in the linear function was 5. For the additive genetic and permanent environmental effects $m = 2$.

Priors. For variance components we considered flat priors : $p(\mathbf{G}), p(\mathbf{P}), p(\sigma_e^2) \propto \text{constant}$. However, for \mathbf{b} we considered alternative priors because previous experience with fixed effects models indicated that the posterior probability of the indicator variables is sensitive to the choice of the priors for \mathbf{b} (Steibel and Grignola, 2000 ; 2001). These priors were the mean and variance of each regressor obtained from a previous run considering the full model. The resulting variance of each regressor was multiplied by a constant, $k = 1, 10, 100, n_i$ (n_i :

number of observations corresponding to each regressor) to vary the amount of information in the prior.

Gibbs sampling. A blocked Gibbs sampling algorithm was used to estimate the linear parameters, and indicator variables were sampled one at a time. The algorithms of Raftery and Lewis (1992) were used to assess convergence. In order to estimate the quantile 0.5 of each parameter with a precision of 0.01, a burn-in period of 500 iterations and a sample size of 50000 iterations were required.

RESULTS AND DISCUSSION

Posterior probabilities. The posterior probabilities of the selected models were affected by the choice of the prior for **b** (figure 1). Based on the variance components estimates and previous experience with this method (Steibel and Grignola, 2001), only results from the analyses using $k = 10$ and $k = n_i$ for the priors of **b** are presented.

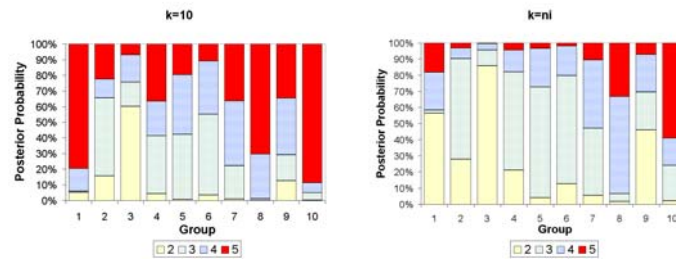


Figure 1. Posterior probabilities for each possible model using different priors for **b. Different colors represent the order of the selected polynomial.**

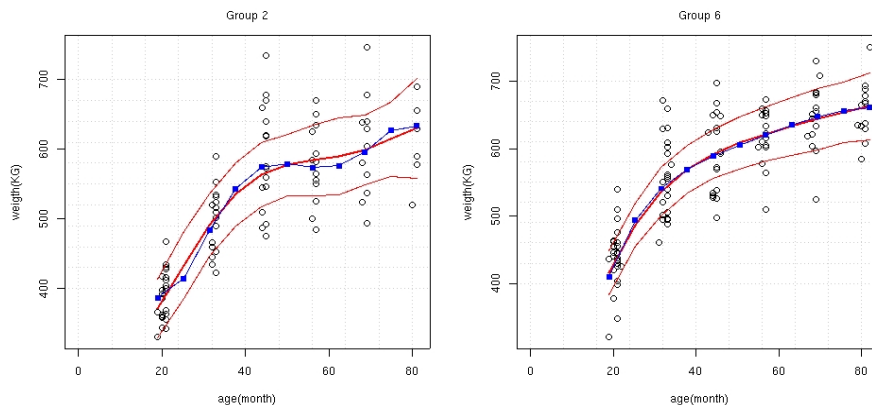


Figure 2. Growth curves posterior means (thick lines) \pm 3 STD (thin lines) obtained across models and $k = 10$, and posterior means for the full model (blue lines with squares)

Average growth curves. Figure 2 shows the posterior probabilities for growth curves averaged across models, for 2 contemporary groups and prior $k = 10$, and for the complete

model. The results indicate that the curve from the average model performs as good as the full model.

Genetic parameters. Quantiles for the posterior distribution of the variance components presented in table 1 show that the Bayesian selection method gives similar results to those obtained under the full model.

Table 1. Quantiles for the marginal posterior distribution of the variance components corresponding to the full models and models with variable selection

Element	Full-model						Selection k=10					
	G			P			G			P		
	Q _{0.05}	Q _{0.50}	Q _{0.95}	Q _{0.05}	Q _{0.50}	Q _{0.95}	Q _{0.05}	Q _{0.50}	Q _{0.95}	Q _{0.05}	Q _{0.50}	Q _{0.95}
11	2145	3392	4609	482	1216	2279	2144	3371	4562	472	1191	2236
12	265	547	837	-17	158	400	263	541	826	-13	157	392
13	-764	-501	-271	-416	-194	-32	-751	-492	-265	-414	-196	-35
22	68	147	245	24	73	159	65	142	237	23	70	153
23	-119	-59	-1	-52	-1	45	-117	-60	-3	-53	-4	41
33	58	116	200	24	69	140	55	111	193	23	69	138
σ_{ϵ}^2	1262	1359	1465	1262	1359	1465	1284	1383	1493	1284	1383	1493

CONCLUSION

The Bayesian variable selection algorithm presented here incorporates the uncertainty about the specification of the linear functions in the inference of the model parameters. It can be considered as an alternative to the Bayesian model averaging proposed by Raftery *et al.* (1997) when the marginal posterior distribution of the parameters are averaged across models. Next steps will include the implementation of the method for the additive genetic and permanent environmental effects.

REFERENCES

- Jensen, J. (2001) *J. Dairy Sci.* 84 : 2803-2812.
 Kuo, L. and Mallick, B. (1998) *Sankhya*, B. **60** : 65-81.
 Meyer, K. (1998) "DFREML User notes".
 Raftery, A., Madigan, D. and Hoeting, J. (1997) *JASA* **92** : 179-191.
 Raftery, A. and Lewis, S.M. (1992) In "Bayesian Statistics 4". Editors Bernardo J.M. Smith A.F.M. Dawid, A.P. and Berger, J.O. Oxford University Press.
 Steibel, J. and Grignola, F. (2000) XVI Reunión Latinoamericana de Producción Animal, III Congreso Uruguayo de Producción Animal. Marzo 28-31, Montevideo, Uruguay.
 Steibel, J.P. and Grignola, F. (2001) *Proc. Reunión Científica Grupo Argentino de Biometría*. p14. Argentina.